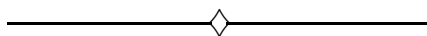


UNIVERSITÄT WÜRZBURG



INSTITUT FÜR THEORETISCHE PHYSIK

AM HUBLAND, D-97074 WÜRZBURG, GERMANY



Phase Transitions of Neural Networks

Wolfgang Kinzel

*Plenary talk for the MINERVA workshop on mesoscopics,
fractals and neural networks, Eilat, March 1997*

Phase transitions of neural networks

Wolfgang Kinzel

Institut für Theoretische Physik

Universität Würzburg, Am Hubland

D-97074 Würzburg, Germany

Abstract

The cooperative behaviour of interacting neurons and synapses is studied using models and methods from statistical physics. The competition between training error and entropy may lead to discontinuous properties of the neural network. This is demonstrated for a few examples: Perceptron, associative memory, learning from examples, generalization, multilayer networks, structure recognition, Bayesian estimate, on-line training, noise estimation and time series generation.

1 Introduction

Since about 15 years there exists a wave of interdisciplinary research activities under the topic "neural networks". Neurobiologists, computer scientists, mathematicians, physicists, psychologists, and linguists are making a more or less common effort to understand the cooperative properties of a system of interacting neurons [Hertz et al 1991]. Meanwhile, the initial excitement and exaggerated promises have been replaced by practical research programs, but much has been achieved and many interesting and unexpected results have been obtained.

The research on neural networks may be classified into three objectives:

1. Neurobiology: The material basis of our brain are about 10^{11} neurons, each of which is directly connected to about 10^3 other ones via synapses. We know a lot about these single units and their anatomical and functional organization. However, we are still far away from understanding learning, association, memory, recognition and generalization on the basis of interacting neurons and their synaptic plasticity. It may be a philosophical problem whether mind, soul, creativity and consciousness can be understood by collective properties of a system of nerve cells. But there is a good chance to elucidate the basic properties of a real neural network by simple models.
2. Computer science: There exists a variety of algorithms which use concepts from real neural networks. Simple units represent information and interact by synaptic weights. Such systems are trained by a set of examples. After the training phase, in which the synaptic weights are adapted to the presented examples, the network is able to achieve a knowledge about the rule which produced the examples; it can generalize. These algorithms are called neural networks or neurocomputer; they are presently applied to a large variety of problems in engineering, science and economy. They have several advantages compared to standard approaches, and there is hope to solve problems by neural networks which are too hard for methods of rule and data based algorithms of artificial intelligence.
3. Physics: Neural networks definitely belong to the class of complex systems, which are characterized by nonlinear dynamics, feedback and macroscopic properties emerging from a huge number of interacting units. In general, physics is interested in understanding such systems in terms of mathematical relations, scaling laws, phase transitions etc.

In physics mathematical modelling of nature has been very successful. How-

ever, it is not clear at all whether such a complex system like the brain can be described by a mathematical language, by simple relations between macroscopic functions and microscopic mechanisms.

On the other side, the full quantum mechanical description of an iron solid is not possible, either. Nevertheless one gains a lot of insight into the spontaneous magnetic ordering below a critical temperature if one studies the Ising model, which replaces the rather complex iron atom by a simple binary unit interacting with its neighbors. With this analogy it is definitely useful to investigate simple units, which model a few essential mechanisms of neurons and synapses, and to study the cooperative behaviour of such interacting units. It is not obvious at all, whether such a system can store an infinite number of patterns with one set of synapses, learn from examples and generalize. Many questions can only be answered from a mathematical calculation.

In this talk I want to emphasize the contribution of statistical physics to the theory of neural computation. Using models and methods from the physics of condensed matter one has been able to calculate the properties of neural networks. This research program uses methods developed already at the beginning of this century by L. Boltzmann and J. W. Gibbs. In 1975 S. F. Edwards and P. W. Anderson, S. Kirkpatrick and D. Sherrington developed a theory of spin glasses [Fischer and Hertz 1991] which was extended to neural networks by J. J. Hopfield [1982]. The first analytic solution of attractor networks succeeded in 1985 [Amit et al 1987]. The statistical mechanics of learning from examples was pioneered by the late E. Gardner [1988]. These approaches opened a new field of research, which produced a lot of interesting results [for reviews see Watkin et al 1993, Oppen and Kinzel 1996]. In view of the big challenge to understand the brain, the statistical physics of neural networks will definitely survive over the next century.

Long before the approach of statistical mechanics, mathematical models of neural networks have been investigated in detail with great success [Hertz

et al 1991]. But in my talk I want to demonstrate on a few examples, to what extent the physics approach is able to ask questions and obtain results which are different from the approach of mathematics and computer science. In particular I think that only the methods of physics can calculate discontinuous and singular properties of infinitely large networks. Hence, in this paper I discuss examples from attractor networks, generalization, structure recognition, Bayesian estimate, on-line training, noise estimation and time series generation, which show discontinuous behavior as a function of model parameters or the size of the training set.

This talk is not supposed to be a review. I apologize to all of my colleagues whose important contributions to the theory of neural networks are not mentioned.

2 Perceptron

The simplest model of a neural network has already been introduced by Rosenblatt in 1960. It consists of an input layer of "neurons" $S_i, i = 1, \dots, N$, which take only binary values $S_i \in \{-1, +1\}$. The activity S_i of each neuron is transmitted by "synapses" $W_i \in \mathbb{R}$ to an output neuron σ as shown in Fig. 1. The output reacts to the sign of the "postsynaptic potential",

$$\sigma = \text{sign} \sum_{i=1}^N W_i S_i = \text{sign } \mathbf{W} \cdot \mathbf{S} \quad (1)$$

In the training phase this network, which was called "perceptron", receives a set of training examples $(\sigma^\nu, \mathbf{S}^\nu)$, $\nu = 1, \dots, \alpha N$. It changes its weight W_i such that a maximum number of examples is correctly mapped by Eq. (1). A simple algorithm has been investigated by Rosenblatt (see Hertz et al 1991): It presents the examples in an arbitrary sequence. If an example is not correctly classified, that is if $\mathbf{W}(t) \cdot \mathbf{S}^\nu \sigma^\nu < 0$, then

$$\mathbf{W}(t+1) = \mathbf{W}(t) + \frac{1}{N} \mathbf{S}^\nu \cdot \sigma^\nu \quad (2)$$

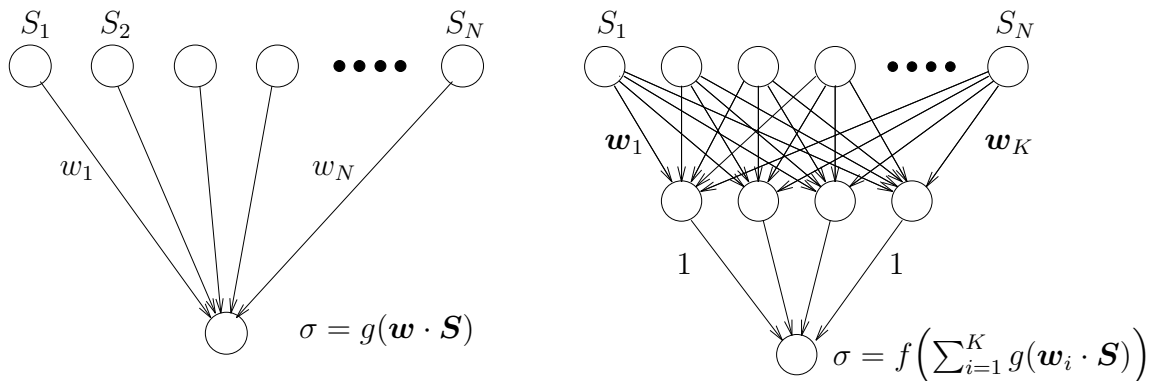


Figure 1: Architecture of the perceptron (left) and the committee machine(right).

There exists a convergence proof for this algorithm: If the examples can be mapped correctly by any perceptron, Eq. (1) with weights \mathbf{W}^* , then the perceptron rule Eq. (2) finds a solution, i. e. the algorithm stops.

The Rosenblatt training rule stems from neurobiology, as proposed by D. Hebb in 1949: Each synapse reacts to the neuronal activities at its two ends. Here we need an additional influence of the postsynaptic potential.

The perceptron implements a linear separable Boolean function, which has an interesting geometrical interpretation: $\mathbf{W} \cdot \mathbf{S} = 0$ defines a hyperplane in an N -dimensional space of inputs \mathbf{S} , the weight vector \mathbf{W} is normal to this plane. On the side of the vector \mathbf{W} the perceptron classifies each input \mathbf{S} to $\sigma = +1$ (black, correct, ...), on the other side the label is $\sigma = -1$ (white, wrong, ...) see Fig. 2.

Now we consider a set of αN many points \mathbf{S}^ν in N dimensions. How many sets of labels $\{\sigma^\nu\}$ can be represented by any perceptron? Surprisingly this problem which is important for the theory of neural computation was already solved by the Swiss mathematician Ludwig Schläfli in the last century [Schläfli 1852]. If any subset of N points is linearly independent, then the number C of possible sets of labels $\{\sigma^\nu\}$ is given by

$$C = 2 \sum_{i=0}^{N-1} \binom{\alpha N - 1}{i}. \quad (3)$$

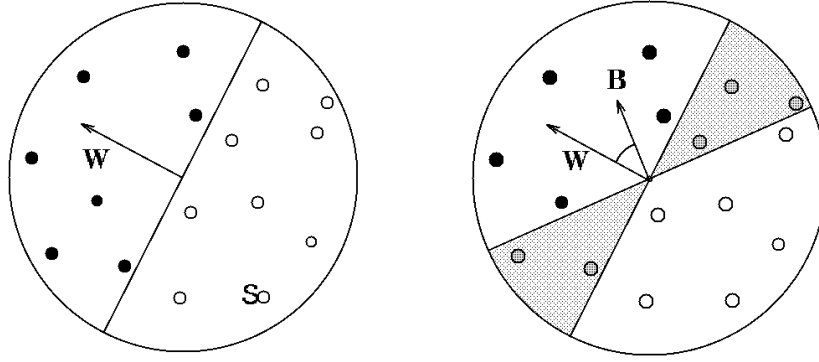


Figure 2: Space of input vectors \mathbf{S} . The weight vector \mathbf{W} of the perceptron defines a hyperplane in the N -dimensional space, which separates the labels σ of the input vectors \mathbf{S} . In the shaded region the labels of teacher \mathbf{B} and student perceptron \mathbf{W} are different.

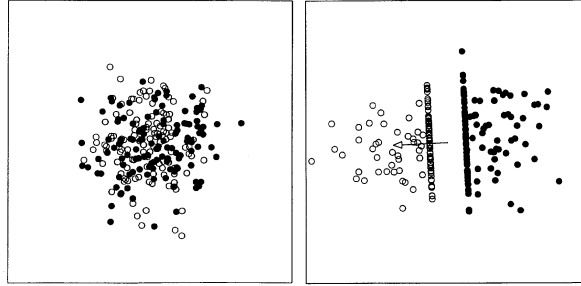


Figure 3: Two dimensional projections of 250 points in 200 dimensions. The points are labelled randomly. The perceptron algorithm finds a hyperplane which separates different labels.

For $\alpha < 1$ all labels can be produced by a perceptron, i. e. $C = 2^{\alpha N}$. For $\alpha < 2$ there is a large fraction of labels which can be separated by a hyperplane. For $\alpha > 2$ only a tiny fraction of patterns can be stored, this fraction disappears for $N \rightarrow \infty$. This result has consequences for the associative memory which will be discussed in the following section: In the limit of $N \rightarrow \infty$ a network with N neurons can store up to $2N$ random patterns.

The geometry of this result is shown in Fig. 3. 250 points are located in

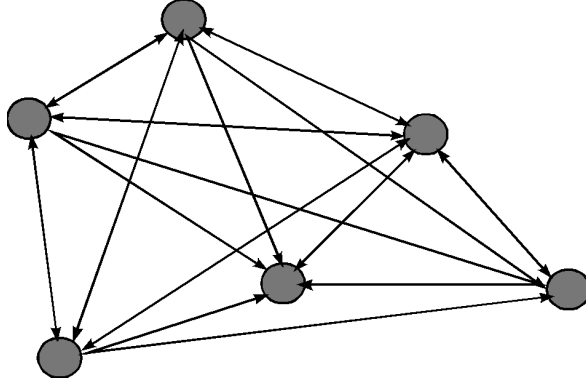


Figure 4: A set of six perceptrons is connected to form an attractor network.

a 200 dimensional space and randomly colored black or white with equal probability. Now we are moving in the space of points and would like to find a view where black and white is clearly separated by a single gap. From Eq. (3) we find with $N = 200$ and $\alpha N = 250$: $C/2^{250} \simeq 1 - 4 \cdot 10^{-23}$; that means for random labels it is almost sure that one can find such a view. In fact the Rosenblatt algorithm, Eq. (2), found the solution shown in Fig. 3.

3 Attractor networks

The perceptron is the "atom" of all neural networks. Many of such elementary units can be composed to a large and complex network. Here we consider an attractor network which consists of N neurons S_i as before. But now each element S_i is connected to any other element S_j by a coupling $W_{ij} \in \mathbb{R}$, as illustrated in Fig. 4.

We want to store αN many patterns $S_i^\nu \in \{-1, +1\}$; $i = 1, \dots, N$; $\nu = 1, \dots, \alpha N$. If we use the Rosenblatt rule, Eq. (2), without the additional condition, we obtain the Hebbian couplings

$$W_{ij} = \frac{1}{N} \sum_{\nu} S_i^\nu S_j^\nu \quad (i \neq j). \quad (4)$$

Since each input S_j is output of a perceptron with weights W_{jk} , we can define a dynamics of the configuration of neurons, \mathbf{S} . For instance, for each neuron

S_i we consider the local field

$$h_i = \sum_j W_{ij} S_j(t) \quad (5)$$

where t is a discrete time index. Now we define a stochastic dynamics by the probability P to find neuron S_i in the state $S \in \{+1, -1\}$ in the next time step $t + 1$:

$$P[S_i(t + 1) = S] = \frac{e^{\beta h_i S}}{2 \cosh(\beta h_i S)} \quad (6)$$

β is a parameter which measures the noise level of the dynamics. For $\beta \rightarrow \infty$ we obtain the noiseless deterministic equation

$$S_i(t + 1) = \text{sign} \sum_j W_{ij} S_j(t). \quad (7)$$

This model was introduced by Hopfield [1982]. He noticed that the dynamics of the neurons is nothing else than the usual Monte Carlo procedure to obtain thermal equilibrium. Since the couplings are symmetric, $W_{ij} = W_{ji}$, the stationary state is given by a Boltzmann distribution

$$P(\mathbf{S}) = \exp(-\beta H(\mathbf{S}))/Z \quad (8)$$

with a Hamiltonian

$$H(\mathbf{S}) = -\frac{1}{2} \sum_{i \neq j} W_{ij} S_i S_j \quad (9)$$

This is the main advantage of equilibrium statistical mechanics: The dynamics $\mathbf{S}(t)$ is replaced by a summation over all possible states \mathbf{S} . Instead of solving a system of N strongly coupled nonlinear equations, one has to calculate the partition sum

$$Z = \sum_{\{\mathbf{S}\}} \exp \left[-\frac{\beta}{2} \sum_{i \neq j} W_{ij} S_i S_j \right]. \quad (10)$$

In the thermodynamic limit of infinitely many neurons, $N \rightarrow \infty$, and infinitely many patterns, $\alpha = \text{const.}$, the partition sum Z was solved exactly by Amit et al [1987] using the replica method. There are two main steps in the calculation:

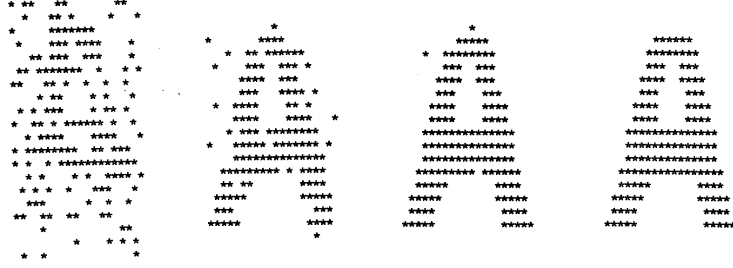


Figure 5: An initial state of an attractor network relaxes to one of the stored patterns. (From Kinzel 1985)

1. The free energy $f = -\ln Z/\beta$ is averaged over all possible sets of patterns $\{\mathbf{S}^\nu\}$. It can be shown that the average value gives the same results as the value f for a single, randomly chosen set of patterns. Hence, this calculation yields results for a typical situation.
2. The sum over 2^N states in $\ln Z$ is performed for fixed order parameters. The minimum of f as a function of these quantities yields their physical values, which describe the stationary state. Hence, the complex system of interacting neurons is described exactly by a few order parameters.

The first step is done using the replica method:

$$\langle \ln Z \rangle_{\{\mathbf{S}^\nu\}} = \lim_{n \rightarrow 0} \frac{\partial}{\partial n} \langle Z^n \rangle_{\{\mathbf{S}^\nu\}}. \quad (11)$$

In the second step one is interested in the overlap between the state \mathbf{S} and one of the patterns \mathbf{S}^ν . Let us assume that the first pattern \mathbf{S}^1 has the form of an "A" as shown in Fig. 4 with $N = 400$. The other 31 patterns consist of random bits. If the initial state $\mathbf{S}(0)$ has an overlap to the first pattern, for instance if it is the noisy "A" of Fig. 5, then after a few steps given by Eq. (7) the network relaxes to the stored information more or less completely.

The statistical mechanics gives information about the possible final states of the dynamics. Here we are interested in the overlap after the relaxation:

$$m_A = \frac{1}{N} \mathbf{S} \cdot \mathbf{S}^A. \quad (12)$$

It turns out that one obtains this order parameter if one calculates the free energy. The overlaps

$$m_\nu = \frac{1}{N} \mathbf{S} \cdot \mathbf{S}^\nu$$

to the other 31 patterns are of the order of $1/\sqrt{N}$. However, their sum r is an additional order parameter

$$r = \frac{1}{\alpha} \sum_{\nu=2}^{\alpha N} m_\nu^2$$

r measures the fluctuation of the final state to the rest of the patterns. Finally there is an order parameter q which measures the complexity of the space of possible solutions \mathbf{S} . Like in the theory of spin glasses, it signals an additional order of the stationary states which has no simple interpretation [Fischer and Hertz 1991].

The theory of attractor networks has close similarities to the theory of an Ising ferromagnet with infinite range interactions. In both cases energy and entropy can be expressed in terms of order parameters. For the ferromagnet one obtains an implicit equation for the spontaneous magnetic order m [Yeomans 1992]

$$m = \tanh \beta J m . \quad (13)$$

For the attractor network one finds

$$m_A = \langle \tanh(\beta m_A + \beta \sqrt{\alpha r} z) \rangle_z \quad (14)$$

where the average is performed over a Gaussian distributed quantity z . There are additional equations for r and q . Hence, compared to the ferromagnet, the $\alpha N - 1$ patterns add a noise term to the local fields.

Fig. 6 shows the result of the analytic calculation [Amit et al 1987]. In the noise-load plane one obtains several phases, which are well separated in the thermodynamic limit. For strong noise, $T = 1/\beta > 1$, or for high load, $\alpha > 0.14$, the network cannot recognize the stored patterns at all. Nevertheless there is a spin glass order for low noise $T < T_g$ with $q > 0$. Only for $T < T_m(\alpha)$ the network can relax to final states which have an

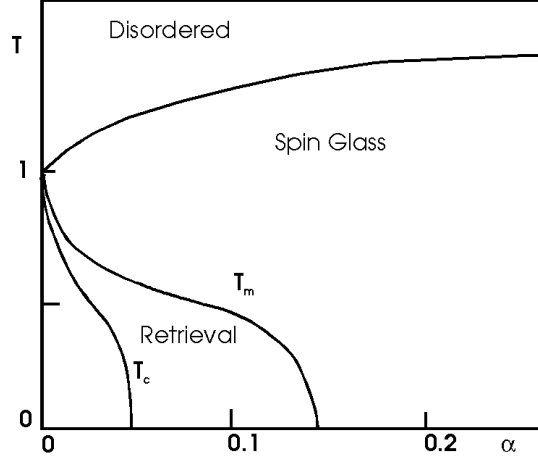


Figure 6: Schematic phase diagram of the Hopfield model. (From Amit et al, 1987)

overlap to one of the stored patterns. This overlap jumps discontinuously to zero at T_m . For $T < T_c(\alpha)$ this retrieval state has the lowest free energy, i. e. it is thermodynamically stable. Note that also for the deterministic dynamics ($T = 0$) there is a discontinuous drop to zero retrieval at $\alpha = \alpha_c \simeq 0.14$. For $\alpha < \alpha_c$ the network restores stored information very well, for $\alpha > \alpha_c$ the network relaxes to final states which have nothing in common with the stored patterns.

According to Schläfli there are couplings with a storage capacity of $\alpha_c = 2$, but these interactions are not symmetric and one cannot apply statistical mechanics to solve the corresponding attractor network. A network with $\alpha_c = 1$ has been analysed by Kanter and Sompolinsky [1987].

In summary, the attractor network functions as an associative memory. It is a distributed memory, since all patterns are stored in all couplings. It is content addressable, since a state with a partial information relaxes to the complete one. Even with a stochastic dynamics it performs well, if the noise level T and storage capacity α are not too high. There is a sharp, discontinuous transition between good and zero performance. Processing of information emerges as a cooperative effect from a large number of simple interacting units.

4 Generalization

We have seen how a network of mutually interacting perceptrons can work as an associative memory. But already the simple perceptron itself has interesting properties. It can learn from examples and recognize an unknown rule.

Consider a perceptron with a weight vector \mathbf{W} , as in Eq. (1). We will call this perceptron the "student". It obtains a set of examples $(\sigma_B^\nu, \mathbf{S}^\nu)$, $\nu = 1, \dots, \alpha N$ from a "teacher". In the simplest case the teacher is another perceptron with a weight vector \mathbf{B} . To what extent can the student gain information about the vector \mathbf{B} if the only available information is the set of αN many examples? The patterns \mathbf{S}^ν are selected randomly and σ_B^ν is the output of the teacher,

$$\sigma_B^\nu = \text{sign } \mathbf{B} \cdot \mathbf{S}^\nu \quad (15)$$

As before we are interested in the limit $N \rightarrow \infty$ and $\alpha = \text{constant}$.

We have to consider two processes:

1. The training phase:

The student network is trained by use of the examples, it tries to decrease the training error

$$\varepsilon_t(\mathbf{W}) = \sum_{\nu=1}^{\alpha N} \Theta[-\sigma^\nu \cdot \sigma_B^\nu] \quad (16)$$

$\Theta(x)$ is the step function, it is zero if the student reproduces the example \mathbf{S}^ν correctly.

2. The test phase:

Now the student receives an input \mathbf{S} which has not been presented before. It gives the answer $\sigma = \text{sign } \mathbf{W} \cdot \mathbf{S}$, which may be different from the answer by the teacher, $\sigma_B = \text{sign } \mathbf{B} \cdot \mathbf{S}$. The probability of disagreement or the generalization error is defined by an average over

all possible input vectors \mathbf{S} :

$$\varepsilon_g = \langle \Theta[-\sigma \sigma_B] \rangle_{\mathbf{S}} \quad (17)$$

From Fig. 2 one can see that ε_g is determined by the angle between the weight vectors of the student and the teacher perceptron

$$\varepsilon_g = \frac{1}{\pi} \arccos \frac{\mathbf{B} \cdot \mathbf{W}}{|\mathbf{B}| |\mathbf{W}|} \quad (18)$$

Training and generalization of the perceptron has been studied in detail using methods of statistical mechanics [see e.g. Watkin et al 1993, Oppen and Kinzel 1996] on a simple scenario where the weights are restricted to binary values, $W_i \in \{+1, -1\}$ and $B_i \in \{+1, -1\}$. The student perceptron is trained by a stochastic algorithm, for instance by a Monte Carlo procedure similar to Eq. (6). But now we have a stochastic dynamics of the synaptic weights W_i instead of the neurons S_i , which leads to a thermal equilibrium given by

$$P(\mathbf{W}) = \exp[-\beta \varepsilon_t(\mathbf{W})]/Z \quad (19)$$

As before we do not have to solve the complex nonlinear dynamics of the weights $\mathbf{W}(t)$ but rather calculate the partition sum

$$Z = \sum_{\{\mathbf{W}\}} \exp[-\beta \varepsilon_t(\mathbf{W})] \quad (20)$$

Again we have to perform two steps:

1. Average $\ln Z$ over all possible sets of examples $\{\mathbf{S}''\}$.
2. Evaluate the sum of 2^N states \mathbf{W} by introducing order parameters.

In the limit of large noise, $\beta \rightarrow 0$, the calculation turns out to be very easy [Seung et al 1992]: The only order parameter is the overlap R between student and teacher,

$$R = \frac{1}{N} \mathbf{B} \cdot \mathbf{W} \quad (21)$$

The training error of Eq. (20) can be replaced by the generalization error

$$\alpha \varepsilon_g = \frac{\alpha}{\pi} \arccos R \quad (22)$$

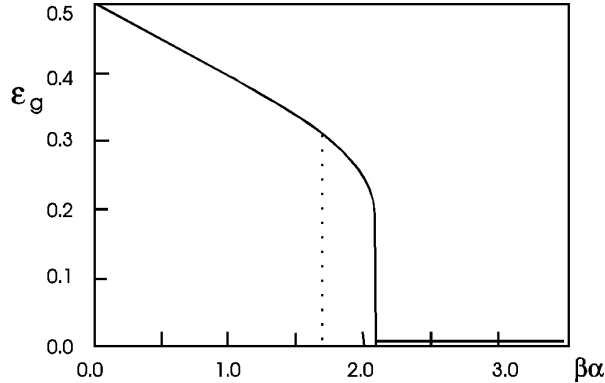


Figure 7: Generalization error ε_g as a function of the size α of the training set (schematic). The binary perceptron is trained stochastically for large noise ($\beta \rightarrow 0$). The dotted line describes the first order phase transition to perfect generalization in thermal equilibrium. The solid line extends to a metastable state. (From Seung et al 1992)

and the entropy is the well known mixture entropy of binary variables

$$S(R) = \frac{1}{2}[(1+R) \ln(1+R) + (1-R) \ln(1-R)] + \ln 2 \quad (23)$$

R is determined by the minimum of the free energy

$$f(R) = \alpha \varepsilon_g(R) - TS(R) \quad (24)$$

Note that the product $\beta\alpha$ in the limit $\beta \rightarrow 0$ is the only parameter of the model, hence a large noise has to be compensated by a large number of examples.

One minimum of f is always $R = 1$, i. e. , the student perfectly recognizes the teacher. However, for $\beta\alpha < 2.08$, the system has an additional minimum at $R < 1$ which is the global one for $\beta\alpha < 1.69$. Fig. 7 shows the generalization error as a function of the fraction of learned patterns. There is a discontinuous transition from poor to perfect generalization, similar to a first order phase transition in solid state physics. Both of the transitions are characterized by metastable states and hysteresis loops. This process of sudden recognition appears even for a noisy training algorithm. A replica calculation shows that the transition qualitatively extends to zero noise $T = 0$. [Seung et al, 1992]

5 Multilayer networks

We have already seen how an attractor network can be built from many perceptrons. Another interesting system which can be constructed from many elementary units is a multilayer network, shown in Fig. 1. It consists of several layers of synaptic weights which map the information coded in the neurons from top to bottom. It is important that such networks can realize any function, if the number of hidden units (neurons in the interior layers) is large enough.

The simplest multilayer network is a committee machine. It consists of N input units, K hidden units, K weight vectors $\mathbf{W}_i, i = 1, \dots, K$ and one output unit σ . The weights of the second layer have the value $+1$, that means, that the output bit σ is given by the majority of the K perceptrons (= members of the committee) in the first layer,

$$\sigma = \text{sign} \left[\sum_{i=1}^K \text{sign } \mathbf{W}_i \mathbf{S} \right] \quad (25)$$

This network is trained from a set of examples $(\sigma_B^\nu, \mathbf{S}^\nu), \nu = 1, \dots, \alpha N$. Note that the opinion of the majority is trained, not the opinion of each member of the committee!

Here we consider the case, where the student is a committee machine with $K = 3$ members. The teacher is a simple perceptron with single weight vector \mathbf{B} . All the weights are assumed to be binary, $W_{ik}, B_i \in \{+1, -1\}$. To what extent can a complex network gain information about a simple rule from a set of examples?

In analogy to the previous section we consider a stochastic training algorithm. The training error ε_t is the "energy" of the Gibbsian probability which describes the stationary state of the stochastic algorithm. In the limit of high noise there are several order parameters, which determine the energy, entropy and generalization error. Firstly, there are the overlaps between the

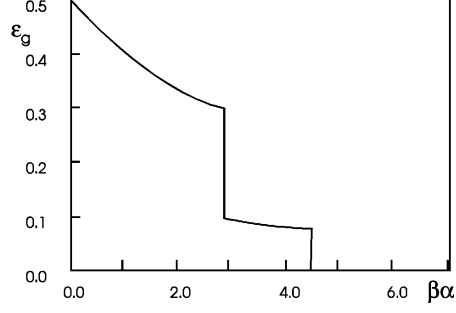


Figure 8: Generalization error ε_g as a function of the size of the training set (schematic). A committee machine with binary weights and three hidden units is trained to a set of examples given by a binary perceptron. (From Schwarze et al 1992)

members of the committee and the teacher,

$$R_i = \frac{1}{N} \mathbf{W}_i \cdot \mathbf{B} \quad (i = 1, 2, 3) \quad . \quad (26)$$

Secondly, the weight vectors of the committee have a mutual overlap

$$Q_{ij} = \frac{1}{N} \mathbf{W}_i \mathbf{W}_j \quad (27)$$

Their physical values are determined by minimizing the corresponding free energy.

Fig. 8 shows the result of the analytic calculation [Schwarze et al 1992]. For a small size of the training set ($\beta\alpha$ small) the members of the committee react symmetrically, $R_1 = R_2 = R_3 < 1$ and the generalization error decreases continuously with α . By increasing the size of the training set, suddenly one of the members recognizes the teacher perfectly, $R_1 = 1, R_2 = R_3 < 1$, and the error jumps to a low value. Further increase of $\beta\alpha$ leads to another discontinuous transition to perfect recognition of the majority. Since the majority vote is already determined by two members, the highest entropy is achieved for $R_1 = R_2 = 1$ and $R_3 = 0$.

Here again the competition between energy and entropy leads to an interesting discontinuous behavior of the generalization ability. Such sharp transitions, which occur for infinitely large networks, only ($N \rightarrow \infty$), are not

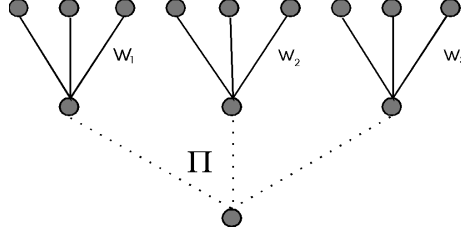


Figure 9: Parity machine with a tree architecture. Each of the three weight vectors \mathbf{W}_i is connected to only one third of the input vector. The output bit is the product of the three hidden units.

obvious. One needs the tools of statistical mechanics to find and describe them.

6 Parity machine

Now we want to discuss another multilayer network, the parity machine with a tree architecture shown in Fig. 9. It consists of N input and K hidden units. The input units are grouped into K parts with N/K neurons each. Each part is input of a perceptron with weights \mathbf{W}_i ($i = 1, \dots, K$). The output of the whole network is given by the parity of the outputs of the K perceptrons,

$$\sigma = \prod_{j=1}^K \text{sign}(\mathbf{W}_j \cdot \mathbf{S}) \quad (28)$$

We consider the case, where both of the student and teacher networks are a parity machine with the same number of units. The teacher network is presenting a set of examples given by

$$\sigma_B^\nu = \prod_{j=1}^K \text{sign}(\mathbf{B}_j \cdot \mathbf{S}^\nu) \quad (\nu = 1, \dots, \alpha N) \quad (29)$$

The examples should be learned without errors. In this case we are interested in the volume V of all student vectors $\{\mathbf{W}_1, \dots, \mathbf{W}_K\}$ which learn the set of examples $\{(\sigma_B^\nu, \mathbf{S}^\nu)\}$ perfectly. V is an integral over a N -dimensional space

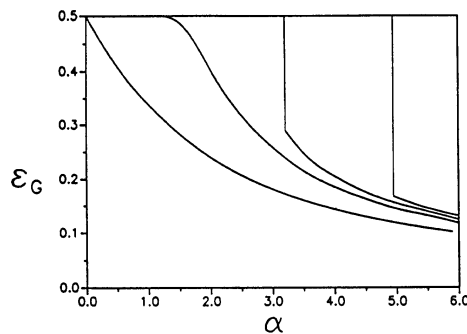


Figure 10: Generalization error as a function of the size of the training set for the parity machine, which learns perfectly a set of examples given by a teacher parity machine. From left to right: $K = 1, 2, 3$ and 4 . (From Oppen 1994)

and corresponds to the partition sum Z of the previous sections. The method of the calculation is similar as before:

1. Average V over all possible sets of inputs $\{\mathbf{S}^\nu\}$ using the replica method.
2. Calculate the integrals by introducing order parameters R and Q , similar to Eqs. (26) and (27).

In general, an additional average of V over all possible teacher vectors \mathbf{B}_i is to be performed.

The generalization error is shown as a function of the number of learned examples in Fig. 10 [Oppen 1994]. It reveals unexpected properties of the network: For a large fraction of examples, $0 < \alpha < \alpha_c(K)$, the network cannot generalize at all ($\varepsilon_g = 1/2$), although it stores of the order of N examples perfectly! Zero training error does not imply an overlap between student and teacher network, even for $\alpha > 0$.

If the number of examples is increased to a critical threshold $\alpha_c(K)$, then the student suddenly recognizes the rule, ε_g jumps to a low value and decreases asymptotically as $1/\alpha$, independently of the number K of hidden units. This property is another surprise: The asymptotic behavior is not determined by

the Vapnik–Chervonenkis dimension, which diverges as $\ln K$ [Barkai et al 1990].

7 Structure recognition

Up to now we have discussed supervised learning, that means the input patterns \mathbf{S}^ν have the labels σ^ν . But there are many applications of neural networks where the labels are not given. In these cases the task is to detect a structure in the data [Hertz et al 1991].

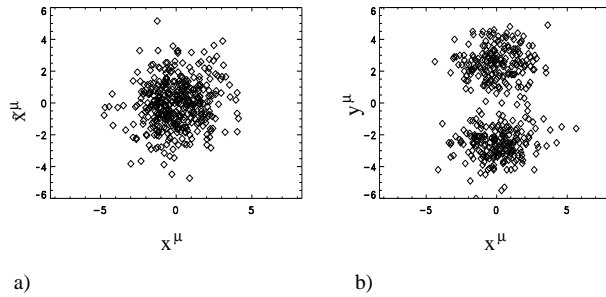


Figure 11: Projection of a distribution of data points. In one direction \mathbf{B} the overlap has a double peak distribution (b), in all orthogonal directions it is Gaussian distributed (a). (From Biehl, 1997)

Consider for example two clouds of αN data points as shown in Fig. 11. The \mathbf{S}^ν are distributed in a N -dimensional space according to a mixture of two Gaussian distributions [Biehl and Mietzner 1993]. This means, that there is a direction \mathbf{B} in data space where the projections $y_\nu = \mathbf{B} \cdot \mathbf{S}^\nu / \sqrt{N}$ have a double peak distribution. In any direction \mathbf{W} orthogonal to \mathbf{B} the corresponding projections x_ν are Gaussian distributed with a single peak, as shown in Fig. 11. Note that the lengths of all vectors \mathbf{S}^ν , \mathbf{W} and \mathbf{B} are of the order of N , while the overlap $\mathbf{S}^\nu \cdot \mathbf{B}$ is of the order of \sqrt{N} .

Given the αN many data points, we want to find the direction \mathbf{B} . There exists a method, well known in engineering, which is called "Principal component analysis" and determines the direction of maximal variance in data

space [Hertz et al 1991]. In fact there is an algorithm for neural networks which finds this direction [Oja 1982]. For our example this means, that we want to find a direction \mathbf{W} which minimizes

$$H(\mathbf{W}) = - \sum_{\nu=1}^{\alpha N} (\mathbf{W} \cdot \mathbf{S}^\nu)^2 / N \quad (30)$$

The minimum of H can be found by calculating the partition sum

$$Z = \int d^N \mathbf{W} \delta(\mathbf{W}^2 - N) \exp(-\beta H(\mathbf{W})) \quad (31)$$

in the limit of $\beta \rightarrow \infty$. Here we have again replaced the dynamics of the algorithm by a summation over all possibilities. As before we have to average $\ln Z$ over all possible data points \mathbf{S}^ν . The evaluation of the N -dimensional integral in the limit $N \rightarrow \infty$ yields the order parameter

$$R = \frac{1}{N} \mathbf{W} \cdot \mathbf{B} \quad (32)$$

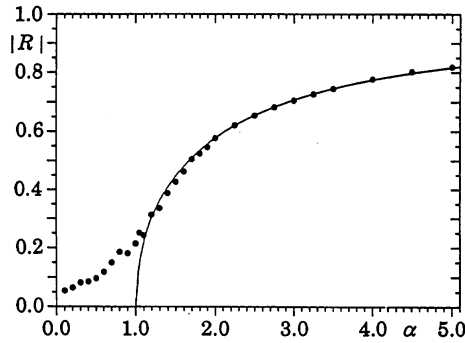


Figure 12: Overlap R between the direction \mathbf{B} of Fig. 11 and the weight vector \mathbf{W} of the training algorithm. The theory (solid line) shows a transition from zero to nonzero recognition with increasing number of data points in the limit $N \rightarrow \infty$. The Monte Carlo simulation (points) for $N = 1000$ show that finite size effects smoothen the transition. (From Biehl and Mietzner, 1993)

The result of this calculation is shown in Fig. 12 [Biehl and Mietzner 1993]. Surprisingly one observes a sharp phase transition. For a small number of data points the system cannot recognize the direction \mathbf{B} of separation at all.

Above a critical number α_c the symmetry $H(\mathbf{W}) = H(-\mathbf{W})$ is spontaneously broken: $|R|$ increases with the deviation from α_c similar to the magnetization in a ferromagnet.

Using the concepts of energy, partition sum and order parameter an unexpected sharp transition from zero to good performance was found for the standard method of principal component analysis.

8 Bayesian estimate

In the previous section we have found the structure of a data distribution by minimizing a cost function. However, if one knows something about the structure of the data it is more efficient to include this knowledge into the algorithm. Here we want to discuss this problem for the two overlapping clouds of data considered in the previous paragraph, see Fig. 11.

The distribution of the data points \mathbf{S}^ν is given by

$$\begin{aligned} P(\mathbf{S}|\mathbf{B};\rho) &\propto \sum_{\tau=\pm 1} \exp\left[-\frac{1}{2}\left(\mathbf{S} - \frac{\rho\tau}{\sqrt{N}}\mathbf{B}\right)^2\right] \\ &\propto \exp[-\beta H(\mathbf{S};\mathbf{B},\rho)] \quad (\beta = 1) \end{aligned} \quad (33)$$

This distribution has two parameters: The vector \mathbf{B} of length N which gives the direction of the cloud separation and the distance ρ between the centers of the clouds.

Now let us assume we know the form of the distribution, Eq.(33), and want to estimate its parameters \mathbf{B} and ρ . Hence our model is for a given distance $\tilde{\rho}$:

$$P(\mathbf{S}|\mathbf{W}) \propto \exp[-\beta H(\mathbf{S};\mathbf{W},\tilde{\rho})] \quad (34)$$

Given the set of data points $\mathbf{S}^\nu, \nu = 1, \dots, \alpha N$, the a posteriori distribution of directions \mathbf{W} is given by the Bayes relation

$$P(\mathbf{W}|\{\mathbf{S}^\nu\}) = \frac{1}{Z} \prod_{\nu=1}^{\alpha N} \delta(\mathbf{W}^2 - N) \exp[-\beta H(\mathbf{S}^\nu; \mathbf{W}, \tilde{\rho})] \quad (35)$$

There are several possibilities to estimate a direction \mathbf{W} from this distribution [Biehl 1997]: Their performance can be measured by the order parameter

$$R = \frac{1}{N} \mathbf{W} \cdot \mathbf{B} \quad (36)$$

which has a single value in the limit $N \rightarrow \infty$. For example one may maximize the a posteriori distribution with respect to \mathbf{W} . The "maximum likelihood" corresponds to the minimum of

$$H(\mathbf{W}) = - \sum_{\nu=1}^{\alpha N} \ln \cosh \frac{\tilde{\rho}}{\sqrt{N}} \mathbf{W} \cdot \mathbf{S}^{\nu} \quad (37)$$

which can be studied by calculating Z for $\beta \rightarrow \infty$, using the replica method [Barkai and Sompolinsky 1994]. For small cloud separation $\tilde{\rho}$ the maximum likelihood solution coincides with first principal component.

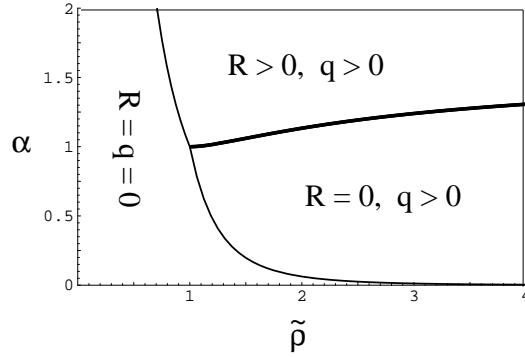


Figure 13: Phase diagram for Bayes estimation of the direction \mathbf{B} of the two clouds of data points in Fig. 11. α measures the number of data points and $\tilde{\rho}$ is the estimated distance of the clouds relative to the true $\rho = 1$. (From Biehl 1997)

Another possibility is to select a direction \mathbf{W} according to the a posteriori distribution, Eq. (35), for instance by using the Monte Carlo method [Watkin and Nadal 1994]. The result is obtained from calculating Z for $\beta = 1$, again by averaging $\ln Z$ over the true distribution of data points. The result of this Gibbs estimate is shown in Fig. 13 [Biehl 1997]. In the plane of α , the size of the data set, and $\tilde{\rho}$, the estimated distance between the clouds, there are

sharp phase transitions. Recognition ($R > 0$) appears only for a sufficiently large number of data points. It is better to use an estimate $\tilde{\rho}$ which is larger than the true one $\rho = 1$, since for small $\tilde{\rho}$ the critical fraction α_c diverges as $\tilde{\rho}^{-4}$.

9 On-line training

In the previous sections **all** of the examples were presented in the training phase of the neural network. The algorithm used the training error with respect to all of the examples in order to find the synaptic weights of the student network. For instance for the Rosenblatt algorithm, Eq. (2) all examples have to be predented several times before the algorithm stops.

Now we want to consider a different training algorithm. At each step only one new example is presented. One does not have to store the complete set of the examples, but the present weight vector \mathbf{W} is changed due to one new example (σ_B', \mathbf{S}') . It turns out that such an "on-line" training is more efficient in terms of computational effort than the "off-line" or "batch" rules of the previous sections, if there are enough examples available.

On-line learning leads to a stochastic differential equation for the weight vector $\mathbf{W}(\nu)$, which becomes a deterministic one for several order parameters in the limit $N \rightarrow \infty$ [Biehl and Schwarze 1995, Saad and Solla 1995]. Usually, the dynamics of on-line learning is not described by a Hamiltonian or a partition sum, nevertheless there are discontinuous properties as a function of model parameters.

Let us consider a two layer network with continuous output neurons. The student as well as the teacher network have three hidden units with weights \mathbf{W}_i and \mathbf{B}_i , the transfer function of the hidden units is the error function.

For simplicity the output neuron is linear with fixed weights

$$\sigma = \sum_{i=1}^3 \text{erf}(\mathbf{W}_i \mathbf{S} / \sqrt{2}) \quad (38)$$

The teacher presents αN many examples given by

$$\sigma_B^\nu = \sum_{i=1}^3 \text{erf}(\mathbf{B}_i \mathbf{S}^\nu / \sqrt{2}) \quad (39)$$

The error of a single example is defined as the quadratic deviation

$$\varepsilon(\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3, \mathbf{S}^\nu) = \frac{1}{2}(\sigma^\nu - \sigma_B^\nu)^2 \quad (40)$$

In analogy to backpropagation [Hertz et al 1991], the change of weights is proportional to the gradient of the training error of a single example:

$$\mathbf{W}_k(\nu + 1) = \mathbf{W}_k(\nu) - \frac{\eta}{N} \vec{\nabla}_{\mathbf{W}_k} \varepsilon \quad (41)$$

From this equation a system of first order, nonlinear and coupled differential equations can be derived for the set of order parameters

$$\begin{aligned} R_{jk} &= \frac{1}{N} \mathbf{W}_j \mathbf{B}_k \\ Q_{jk} &= \frac{1}{N} \mathbf{W}_j \mathbf{W}_k \end{aligned} \quad (42)$$

In our case there are 15 order parameters which change after each presentation of a new example. In the limit of $N \rightarrow \infty$ the index ν becomes a continuous "time" α . Hence, one has to calculate the flow of $R_{jk}(\alpha), Q_{jk}(\alpha)$ in the 15 dimensional space of order parameters which determine the generalization error $\varepsilon_g(\alpha)$.

It turns out that there are several fixed points of this flow, which have important consequences for the behavior of the generalization error. Fig. 14 shows a typical example. For a small number of examples, ε_g decreases fast. But then the generalization error almost does not change for a long training period. Suddenly it decreases to good performance of the network.

This plateau of $\varepsilon_g(\alpha)$ which is observed in applications, too, can be understood in terms of the flow of order parameters [Biehl et al 1997]. There is one

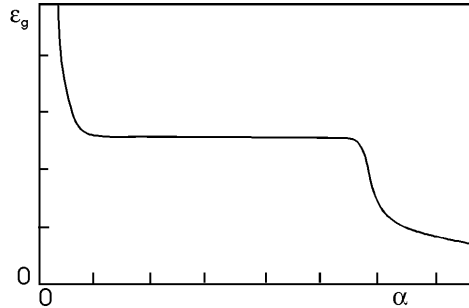


Figure 14: Generalization error as a function of the size of training set for a two layer network (schematic). The plateau is related to a fixed point with one weak repulsive direction for the flow of order parameters. (From Biehl et al 1997)

fixed point which is strongly attractive in almost all directions. But in one or a few directions it has a small repulsive component. Hence, the flow remains for quite a while (depending on initial conditions) close to this fixed point with a large generalization error, but then it flows away to the completely attractive fixed point with zero error.

The number of fixed points depends on the learning rate η of the training rule. In our simple example there are at least ten different fixed points for $\eta = 1$. With increasing η some fixed points split into two, which usually means that some symmetry is broken. Fixed points suddenly appear, disappear or annihilate with varying learning rate η . Such discontinuous behavior is reflected in the generalization error $\varepsilon_g(\alpha)$.

10 Noise estimation

The examples, given by a teacher network, may have errors. To what extent can a student network derive information about the teacher weights from a set of faulty examples? This problem has been investigated in detail [Copelli et al 1997]. We consider a committee machine with a tree architecture

$$\sigma = \text{sign}\left[\sum_{i=1}^K \text{sign } \mathbf{W}_i \mathbf{S}_i\right] \quad (43)$$

The student as well as the teacher network have the same number K of hidden units. The examples are distorted by noise: The bit σ_B^ν is flipped with probability λ . For $\lambda = 1/2$ there is no information in the examples, but for $0 < \lambda < 1/2$ the student network may obtain an overlap to the teacher one with increasing number α of examples.

We study the training algorithm

$$\mathbf{W}_k(\nu + 1) = \mathbf{W}_k(\nu) + \frac{1}{N} F_k \mathbf{S}_k \quad (44)$$

Instead of a parameter η we have introduced a function F_k which is determined from a variational principle which maximizes the decrease of generalization error ε_g at each training step [Kinouchi and Caticha 1992]. Hence, one can define an optimal weight change, which depends on the order parameters. It also contains the noise rate λ which is not known in general; it has to be estimated by a value Λ , $F_k(\lambda)$ is replaced by $F_k(\Lambda)$.

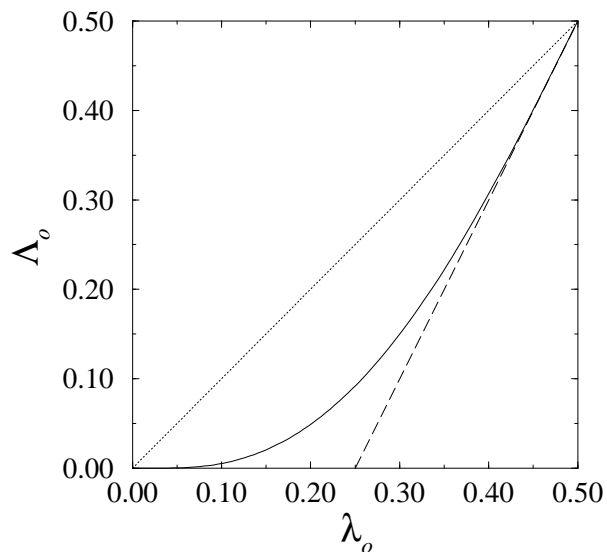


Figure 15: On-line training with optimal weights. Λ is the estimated and λ the true noise level of the training examples. Above the solid line perfect generalization is possible, $\varepsilon_g(\alpha \rightarrow \infty) = 0$. Below the dashed line the network cannot generalize at all. Between the dashed and solid line a partial recognition of the teacher network is possible. (From Copelli et al, 1997)

The generalization error has been calculated for optimal on-line learning [Copelli et al 1997], Fig. 15 shows the result. There are sharp boundaries in the (λ, Λ) plane where the behavior of the network changes drastically. Fastest decay of $\varepsilon_g(\alpha)$ is obtained if the true noise is estimated correctly, $\lambda = \Lambda$, as expected. If the estimated noise parameter λ is large enough, then $\varepsilon_g(\alpha)$ still decreases to zero in the limit of an infinite number α of examples. However, if Λ is below the dashed line, then the network cannot generalize at all. If Λ lies in the intermediate region then the generalization error decreases to a nonzero value; the network can generalize only partially. Again we observe sudden changes in the behavior of the network as a function of model parameters.

11 Time series generalization

Most of the work on the statistical physics of neural networks has been done on static data. A set of input vectors $\{\mathbf{S}^\nu\}$ is taken from a static distribution and classification labels $\{\sigma^\nu\}$ are taken from a static rule. Only recently this research program has been extended to the analysis of time series [Eisenstein et al, 1995], which is an important field of applications of neural networks.

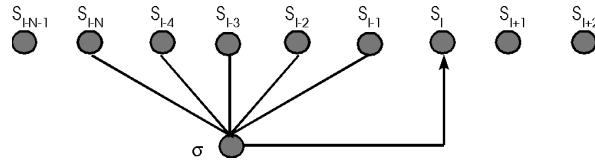


Figure 16: A perceptron working as sequence generator. The sequence S_i of numbers is generated by a perceptron moving to the right.

In the simplest case our elementary unit, the perceptron, is trained to a sequence S_i of real numbers, where i is a discrete time. As shown in Fig. 16, the perceptron takes a window of N numbers as input and makes a prediction

of the following number,

$$S'_l = \tanh\left[\frac{\beta}{N} \sum_{j=1}^N W_j S_{l-j}\right] \quad (45)$$

In the training phase the weights \mathbf{W} are changed to decrease the error $(S'_l - S_l)^2$.

In order to apply the concepts of statistical physics one needs a well defined sequence $\{S_i\}$. As before it may be given by a teacher perceptron with weight vector \mathbf{B} . For a given window of N numbers $(S_{l-N}, \dots, S_{l-1})$ the perceptron defines the number S_l . Then it moves one step and generates S_{l+1} . It turns out that the generation of time series is already an interesting problem with many unsolved puzzles [Eisenstein et al 1995, Kanter et al 1995, Schröder and Kinzel 1997].

The numerical investigation of the sequence generator shows that an initial state of Eq (45) approaches a quasi periodic attractor which is related to a peak in the Fourier spectrum of \mathbf{W} [Eisenstein et al 1995]. Hence, the perceptron selects one mode of the couplings. Therefore it is useful to study couplings with a single Fourier component

$$W_j = \cos(2\pi k \frac{j}{N} - \pi\phi) \quad (46)$$

k is the frequency and ϕ the phase of the weights. An attractor of Eq. (45) is the solution of

$$S_l = \tanh\left[\frac{\beta}{N} \sum_{j=1}^N \cos(2\pi k \frac{j}{N} - \pi\phi) S_{l-j}\right] \quad (47)$$

Recently this equation could be solved analytically [Kanter et al 1995]: For small values of β the attractor is zero, $S_l = 0$. For a critical value, which is independent of the frequency k ,

$$\beta_c = 2 \frac{\pi\phi}{\sin \pi\phi} \quad (48)$$

there appears the solution

$$S_l = \tanh\left[A(\beta) \cos(2\pi(k + \phi)\frac{l}{N})\right] \quad (49)$$

The phase ϕ of the weights shifts the frequency $k + \phi$ of the solution. The amplitude $A(\beta)$ is a continuous function of the distance $\beta - \beta_c > 0$ to the critical point. For $\beta \rightarrow \infty$ \tanh is replaced by sign and the sequence generator becomes a bit generator. In this case the solution, Eq. (49), is more complicated, but again the phase ϕ shifts the frequency of the bit sequence [Schröder and Kinzel 1997].

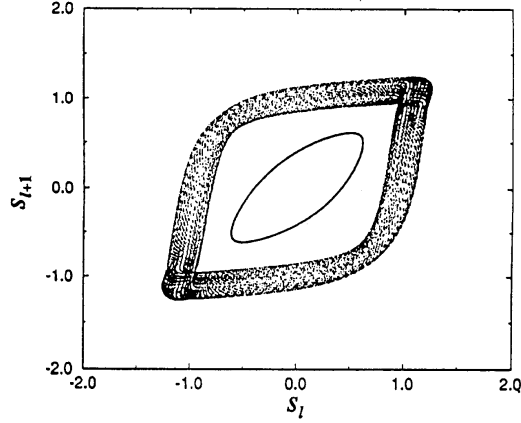


Figure 17: Return map of a sequence generated by a multilayer network with two hidden units. The one-dimensional attractor in the middle becomes a two-dimensional one, if the parameter β is increased. (From Kanter et al, 1995)

If the network is a multilayer perceptron with K hidden units, then the attractor of the sequence generator is a kind of superposition of the single modes of each hidden unit [Kanter et al 1995]. Each unit has its own critical point and the number of nonzero solutions determines the dimension of the attractor. This is shown in the return map of Fig. 17 where S_{t+1} is plotted against S_t for $K = 2$ hidden units. Increasing β first one component is activated leading to a one-dimensional attractor shown in the middle of the figure. Since in general the frequency is irrational, the attractor is a continuous curve. For larger values of β the second unit is activated, giving the two-dimensional attractor with larger amplitude.

12 Summary

Concepts of statistical physics have successfully been applied to the theory of neural computation. The cooperative behavior of a large number of interacting neurons can be described in terms of partition sums and order parameters. The competition between training error and entropy may lead to discontinuous properties.

The approach of statistical mechanics has several advantages:

1. Networks with an infinite number of neurons and synapses can be calculated analytically. Complex cooperative behavior of interacting neurons is described in terms of a few order parameters.
2. The results are obtained for a typical situation, for instance for the most general set of examples.
3. One obtains exact mathematical relations between the observed cooperative properties of the network, its model parameters and the size of the training set.
4. Many networks and algorithms show discontinuous properties as a function of model parameters or the number of presented examples. Statistical physics can describe such sudden changes in the cooperative behavior of the network.

Statistical mechanics of neural networks has been applied to a variety of problems; we just want to mention learning from examples, generalization, associative memory, attractor networks, structure recognition, clustering, classification, coding and time series analysis. For all of these problems general properties have been calculated mathematically. Novel and unexpected results have been found. Hence, I think that in the last 15 years theoretical physics has successfully contributed to our understanding of neural networks, with impact on neurobiology and computer science.

Acknowledgement: The author thanks Michael Biehl for comments on the manuscript.

Literature

Hertz, J. A., A. Krogh and R. G. Palmer, 1991, Introduction to the Theory of Neural Computation (Addison Wesley)

Amit, D. , H. Gutfreund and H. Sompolinsky, 1987, Annals of Physics, **173**, 30-67

Barkai, E. , D. Hansel and I. Kanter, 1990, Phys. Rev. Lett. **65** 2312

Barkai, N. and H. Sompolinsky, 1994, Phys. Rev. **E 50**, 1766

Biehl, M. and A. Mietzner 1993, Europhys. Lett. **24**, 421–426

Biehl, M. and H. Schwarze, 1995, J. Phys. **A 28**, 643(?)

Biehl, M. , 1997, lecture notes, University of Würzburg

Biehl, M. , P. Riegler and C. Wöhler, 1997, to be published

Copelli, M. , R. Eichhorn, O. Kinouchi, M. Biehl, R. Simonetti, P. Riegler and N. Caticha, 1997, Europhys. Lett. **37**, 427-432

Eisenstein, E. , I. Kanter, D. Kessler and W. Kinzel, 1995, Phys. Rev. Lett. **74**, 6

Fischer, K. H. and J. A. Hertz, 1991, Spin Glasses, (Cambridge University Press)

Gardner, E. , 1988, J. Phys. **A 21**, 257–270

Hopfield, J. J. , 1982, Proceedings of the National Academy of Sciences, USA **79**, 2554–2558

Kanter, I. , D. A. Kessler, A. Priel and E. Eisenstein, 1995, Phys. Rev. Lett. **75**, 2614–2617

- Kanter, I. and H. Sompolinsky, 1987, Phys. Rev. **A 35**, 380
- Kinzel, W. , 1985, Z. Physik **B60**, 205
- Kinouchi, O. and N. Caticha, 1992, J. Phys. **A 25**, 6243
- Oja, E. , 1982, J. Math. Biol. **15**, 267
- Opper, M. , 1994, Phys. Rev. Lett. **72**, 2113
- Opper, M. and W. Kinzel, 1996, in Models of Neural Networks III, ed. by E. Domany, J. L. van Hemmen and K. Schulten (Springer, Berlin)
- Saad, D. and S. A. Solla, 1995, Phys. Rev. Lett. **74**, 4337
- Schläfli, L. , 1852, Theorie der vielfachen Kontinuität, Gesammelte Mathematische Abhandlungen, ed. Steiner–Schläfli–Komitee Basel, Birkhäuser p 171
- Schröder, M. and W. Kinzel, 1997, to be published
- Schwarze, H. , M. Opper and W. Kinzel, 1992, Phys. Rev. **A 46**, R 6185
- Seung, H. S. , H. Sompolinsky and N. Tishby, 1992, Phys. Rev. **A 45**, 6056
- Watkin, T. L. H. , A. Rau and M. Biehl, 1993, Rev. Mod. Phys. **65**, 499
- Watkin, T. L. H. and J. –P. Nadal, 1994, J. Phys. **A 27**, 1889
- Yeomans, J. , 1992, Statistical mechanics of phase transitions, (Clarendon Press, Oxford)